# CONSISTENCY OF $\ell_1$ -REGULARIZED MAXIMUM-LIKELIHOOD FOR COMPRESSIVE POISSON REGRESSION

Yen-Huan Li and Volkan Cevher

# LIONS, EPFL

# ABSTRACT

We consider Poisson regression with the canonical link function. This regression model is widely used in regression analysis involving count data; one important application in electrical engineering is transmission tomography. In this paper, we establish the variable selection consistency and estimation consistency of the  $\ell_1$ -regularized maximum-likelihood estimator in this regression model, and characterize the asymptotic sample complexity that ensures consistency even under the compressive sensing setting (or the  $n \ll p$  setting in high-dimensional statistics).

*Index Terms*— Poisson regression, maximum likelihood,  $\ell_1$ -regularization, consistency, variable selection, sample complexity, transmission tomography

### 1. INTRODUCTION

The Poisson regression model with the canonical link function is pervasive in regression analysis for count data [5, 15]. Transmission tomography is one famous application of this model in electrical engineering [10], and this is the main motivation of this work. In this paper, we discuss the variable selection consistency and estimation consistency of  $\ell_1$ regularized maximum-likelihood (ML) estimators in this regression model. We expect that our work will validate, in a theoretically sound fashion, the use of the  $\ell_1$ -regularized ML estimator in the Poisson regression model.

Below is a summary of our main contributions.

- 1. We provide *non-asymptotic* performance guarantees for the  $\ell_1$ -regularized ML estimator in the Poisson regression model; that is, we not only prove the variable selection consistency and estimation consistency when the sample size increases to infinity, but also provide explicit bounds on estimation error and probability of correct variable selection for *any finite value* of the sample size.
- 2. We characterize the scaling of (n, p, s) that ensures variable selection consistency and estimation consis-

tency, where *n* denotes the sample size, *p* denotes the parameter dimension, and *s* denotes the number of non-zero entries in the parameter. Our result shows that the  $\ell_1$ -regularized ML estimator is consistent even when the parameter dimension increases exponentially with the sample size (cf. Corollary 4.1).

3. We derive novel inequalities for *self-concordant like* functions. Our framework enables a structured derivation of the consistency results, in the sense that most parts of our proof can be directly extended to statistical models involving a self-concordant like function, such as the logistic regression model [2]. We also develop computationally efficient algorithms to approximate the  $\ell_1$ -regularized ML estimator for transmission tomography with the theory of self-concordant like functions; due to the page limit, we are not able to show the optimization theoretic results in this paper.

# Notations

We only point out some notations that might cause confusions without explicit definitions; other notations should be standard. For a vector  $v \in \mathbb{R}^p$ , we define the support function  $\operatorname{supp}(v) := \{i : v_i \neq 0\}$ . Let S be a subset of  $\{1, \ldots, p\}$ . We denote by  $v_S$  and  $v_{S^c}$  the sub-vector of v with entries indexed by S and  $\{1, \ldots, p\} \setminus S$ , respectively. Similarly, for a matrix  $A \in \mathbb{R}^{p \times p}$ , we define the sub-matrices  $A_{S,S}$ ,  $A_{S^c,S}$ , etc. We write  $\mathbb{R}^S$  for the |S|-dimensional subspace of  $\mathbb{R}^p$  indexed by S. We shall consider  $\ell_p$ -norms  $\|v\|_p := (\sum_{i=1}^p |v_i|^p)^{1/p}$ , and  $\|A\|_p$  denotes the operator norm of A induced by the  $\ell_p$ -norm.

# 2. RELATED WORKS

The proof of our main theorem starts with a generalization of Wainwright's primal-dual witness approach [26] for the Gaussian linear regression model. However, due to the nonlinearity of the Poisson regression model under consideration, we have to bound a residual term which will be defined clearly in Section 5, as other generalizations of Wainwright's work did [21, 22]. Unlike [21, 22], where specific techniques are developed for different statistical models, our proof is

This work was supported in part by the European Commission under Grant MIRG-268398, ERC Future Proof, SNF 200021-132548, SNF 200021-146750 and SNF CRSII2-147633.

*structured*, in the sense that it applies not only to the Poisson model, but can be easily adapted for any statistical model with a self-concordant like negative log-likelihood function.

The notion of self-concordant like functions was first proposed by Bach in [2], but the inequalities for self-concordant like functions presented in Section 6 are new. Our definition is slightly more general than the one in [2], as our definition is valid for any *norm* on  $\mathbb{R}^p$  (cf. Definition 6.1), while the definition in [2] considers only the Euclidean norm. This allows us to obtain an estimation error bound in the  $\ell_{\infty}$ -norm, while the error bound in [2] is in the  $\ell_2$ -norm. Note that a self-concordant like function is not necessarily self-concordant (cf. [18, 19] for the definition of self-concordant functions).

There are some general theoretical frameworks for proving the consistency of  $\ell_1$ -regularized ML estimators, but they are not directly applicable for our purpose. The framework in [17] relies on the notion of *restricted strong convexity (RSC) condition*. However, the RSC condition alone does not guarantee variable selection consistency. A general framework for verifying variable selection consistency in generalized linear models is provided in [8]. The assumptions therein, however, are difficult to check for specific statistical models. We note that the only seemingly intractable assumptions in our paper are the irrepresentability condition and Assumption 1. The irrepresentability condition [26, 28], which has been proved to be almost necessary for the Gaussian linear regression model [26], and is also present in [8]. We shall discuss the validity of Assumption 1 in Section 4.

There are some other papers about regularized Poisson regression [12, 16, 20], but the regression model they considered is different from the one in this paper. Specifically, the negative log-likelihood functions in [12, 16, 20] are self-concordant but not self-concordant like; thus our analysis does not directly apply to their setting, and vice versa. Moreover, the variable selection consistency is not considered in [12, 16, 20].

We note that the optimal choice of the regularization coefficient based on given samples is in general an open problem [4]. Though in this paper we provide an explicit value of the regularization coefficient that ensures consistency, it is impractical since it requires at least the information of the true sparsity level (cf. (1)). In practice, one may apply the covariance penalty approach or cross validation procedures to estimate the optimal value of the regularization coefficient [1, 7].

## 3. PROBLEM FORMULATION

We formulate the Poisson regression model and the associated  $\ell_1$ -regularized ML estimator in this section. Since the Poisson regression model finds applications in a variety of research areas, we deliberately keep the formulation general and ignore some physical constraints such as the positivity of the parameters in transmission tomography problems [10]. Let  $\theta^{\natural} \in \mathbb{R}^p$  be the parameter to be estimated. Let  $a_1, \ldots, a_n$  be given vectors in  $\mathbb{R}^p$ . The measurement outcomes are samples  $y_1, \ldots, y_n$  of independent real-valued Poisson random variables  $Y_1, \ldots, Y_n$ , with probability distribution given by

$$\mathbb{P}\left\{Y_i = y\right\} = \frac{\exp(-\lambda_i)\lambda_i^y}{y!}, \quad \text{for all } y \in \mathbb{N} \cup \left\{0\right\},$$

where  $\lambda_i := \exp(-\langle a_i, \theta^{\natural} \rangle)$ . Our aim is to estimate  $\theta^{\natural}$  given  $a_1, \ldots, a_n$  and  $y_1, \ldots, y_n$ .

The  $\ell_1$ -regularized ML estimator, denoted by  $\hat{\theta}_n$ , is defined as

$$\hat{\theta}_n \in \arg\min_{\theta \in \mathbb{R}^p} \left\{ f_n(\theta) + \rho_n g(\theta) \right\}$$

for some positive regularization coefficient  $\rho_n$ , where  $f_n$  is the normalized negative log-likelihood (up to some constant shift),

$$f_n(\theta) := rac{1}{n} \sum_{i=1}^n \left[ y_i \langle a_i, \theta \rangle + \exp\left(-\langle a_i, \theta \rangle \right) \right],$$

and g is the  $\ell_1$ -norm,  $g(\theta) := \|\theta\|_1$ . The estimator  $\hat{\theta}_n$  exists because  $f_n + \rho_n g$  is coercive; it might not be unique, so we define it via an inclusion relation.

### 4. MAIN RESULT

In this section we discuss the assumptions, and then show the main theorem.

Assumption 1. The restricted Hessian  $\left[\nabla^2 f_n(\theta^{\natural})\right]_{S,S}$  satisfies

$$\left\langle v, \left[\nabla^2 f_n(\theta^{\natural})\right]_{\mathcal{S},\mathcal{S}} v \right\rangle \ge \mu \left\|v\right\|_1 \left\|v\right\|_{\infty}$$

with some  $\mu > 0$ , for any  $v \in \mathbb{R}^s$ .

*Remark.* Note this assumption also implies the positive definiteness of the restricted Hessian. By (3), f restricted on  $\mathbb{R}^{S}$  is strictly convex.

Assumption 1 is similar to the restricted eigenvalue condition [3], the compatibility condition [25], and the RIP-1 condition [11]. While Assumption 1 cannot be implied by any of them, numerical experiments strongly suggest that Assumption 1 might hold with high probability when  $a_1, \ldots, a_n$  are independent vectors of i.i.d. standard Gaussian random variables and  $n \sim s \log(p/s)$ . We leave the verification of Assumption 1 as a future work.

The second assumption is known as the irrepresentability condition [26, 28]. This condition has been proved to be "almost necessary" for the Gaussian linear regression model in [26].

Assumption 2. The Hessian satisfies

$$\left\| \left[ \nabla^2 f_n(\theta^{\natural}) \right]_{\mathcal{S}^c, \mathcal{S}} \left[ \nabla^2 f_n(\theta^{\natural}) \right]_{\mathcal{S}, \mathcal{S}}^{-1} \right\|_{\infty} < 1 - \alpha$$

for some constant  $\alpha \in (0, 1)$ .

The final assumption is standard [2, 8, 26].

Assumption 3. Let  $A \in \mathbb{R}^{n \times p}$  be the matrix whose *i*-th row is given by  $a_i$ . We assume that each column of A has  $\ell_2$ -norm less than  $\sqrt{n}$ .

**Theorem 4.1.** Suppose Assumptions 1–3 are satisfied. If  $\rho$  is chosen such that

$$\rho_n \le \min\left\{\frac{\alpha\mu^2}{2\lambda_{\max}(4+\alpha)^2 s \left\|A_{\mathcal{S}}\right\|_{\infty}}, \frac{2\mu}{(4+\alpha)\left\|A_{\mathcal{S}}\right\|_{\infty}}\right\}.$$
(1)

then with probability at least  $1 - 2p \exp\left[-\frac{\alpha\sqrt{n\rho}}{16\overline{\lambda}}\right]$ ,  $\hat{\theta}$  satisfies  $\hat{\theta}_{S^c} = 0$ , and

$$\left\|\hat{\theta}_n - \theta^{\natural}\right\|_{\infty} \le \varepsilon_n := \left(\frac{4+\alpha}{2\mu}\right)\rho$$

In addition, if  $\theta_{\min} := \min_{i \in S} \left\{ \left| \theta_i^{\natural} \right| \right\} > \varepsilon_n$ , then  $\hat{\theta}_n$  recovers the sign pattern of  $\theta^{\natural}$ .

Consider the high-dimensional setting, where s and p can scale with n [8, 9, 26, 28].

**Definition 4.1.** A sequence of estimators  $\{\hat{\theta}\}_{n \in \mathbb{N}}$  is consistent in variable selection if there exists  $\{\rho_n\}_{n \in \mathbb{N}}$  such that

$$\lim_{n \to \infty} \mathbb{P}\left\{ \operatorname{supp}\left(\hat{\theta}_n\right) \neq \operatorname{supp}\left(\theta^{\natural}\right) \right\} = 0.$$

**Definition 4.2.** A sequence of estimators  $\{\hat{\theta}\}_{n \in \mathbb{N}}$  is consistent in estimation if there exists  $\{\rho_n\}_{n \in \mathbb{N}}$  such that for any  $\epsilon > 0$ ,

$$\lim_{n \to \infty} \mathbb{P}\left\{ \left\| \hat{\theta}_n - \theta^{\natural} \right\| > \epsilon \right\} = 0$$

where  $\|\cdot\|$  is some norm on  $\mathbb{R}^p$ .

Choose  $\rho_n$  such that

$$n^{-1/2}\log p \ll \rho_n \ll s^{-1} \|A_{\mathcal{S}}\|_{\infty}^{-1}.$$

We obtain the following lemma.

**Corollary 4.1.** If  $||A_{\mathcal{S}}||_{\infty} s \log p \ll \sqrt{n}$  and  $\theta_{\min} > \varepsilon_n$ , then  $\{\hat{\theta}\}_{n \in \mathbb{N}}$  is consistent in variable selection. If in addition,  $\rho_n \to 0$  as  $n \to \infty$ , then  $\{\hat{\theta}\}_{n \in \mathbb{N}}$  is also consistent in estimation.

Note that  $\mu$  and  $\alpha$  are assumed to be constants.

#### 5. SKETCH OF THE PROOF

Due to the page limit, we do not show the complete proof but briefly summarize the logical structure here. We omit the subscripts n in this section, and define  $S := \operatorname{supp}(\theta^{\natural})$ , to simplify the notation. It is desirable to have  $\hat{\theta}$  behave like the oracle estimator  $\check{\theta}$ , defined as

$$\check{\theta} := \arg\min_{\theta \in \mathbb{R}^p: \theta_{S^c} = 0} \left\{ f(\theta) + \rho g(\theta) \right\}.$$

Since  $f + \rho g$  restricted to  $\mathbb{R}^{S}$  is still coercive,  $\check{\theta}$  exists. It can be verified that f is self-concordant like (cf. Definition 6.1). By Assumption 1 and (3),  $f + \rho g$  is strictly convex on  $\mathbb{R}^{S}$  and thus  $\check{\theta}$  is uniquely defined.

The following result is obtained by the primal-dual witness approach and can be proved as in [26], so we omit the proof here.

**Lemma 5.1.** We have 
$$\hat{\theta} = \check{\theta}$$
 if  $\| [\nabla f(\check{\theta})]_{S^c} \|_{\infty} < \rho$ .

By a Taylor series expansion on the first-order optimality condition of  $\check{\theta}_n$  and the triangle inequality, the condition in Lemma 5.1 is satisfied if the irrepresentability condition (Assumption 2) holds, and [13]

$$\max\left\{\left\|\nabla f(\theta^{\natural})\right\|_{\infty}, \left\|r\right\|_{\infty}\right\} \leq \frac{\alpha}{4}\rho,$$

where the residual term r satisfies [27],

$$\|r\|_{\infty} \leq \left\|\check{\theta} - \theta^{\natural}\right\|_{\infty} \sup_{t \in [0,1]} \left\{ \left\|\nabla^2 f(\theta_t) - \nabla^2 f(\theta^{\natural})\right\|_{\infty} \right\},\$$

where  $\theta_t := \theta^{\natural} + t(\check{\theta} - \theta^{\natural})$ . For the Gaussian linear regression model considered in [26],  $r \equiv 0$ , while the rest of our proof is mainly devoted to evaluating the non-zero residual term.

The norm of the residual term r is small if  $\mathring{\theta}$  is close to  $\theta^{\natural}$ . This intuition is quantified by the following lemma.

**Lemma 5.2.** We have  $||r||_{\infty} \leq (\alpha/4)\rho$  if

$$\left\|\check{\theta} - \theta^{\natural}\right\|_{\infty} \le \min\left\{\sqrt{\frac{\alpha\rho}{8s\lambda_{\max}\left\|A_{\mathcal{S}}\right\|_{\infty}}}, \frac{1}{\left\|A_{\mathcal{S}}\right\|_{\infty}}\right\}.$$
 (2)

where  $\lambda_{\max} := \max \{\lambda_1, \dots, \lambda_n\}$ , and  $A_S$  denotes the submatrix of A whose columns are indexed by S.

To verify (2) is equivalent to evaluating the estimation error of the oracle estimator  $\check{\theta}$ . In fact  $\|\check{\theta} - \theta^{\natural}\|_{\infty}$  corresponds to the estimation error of the  $\ell_1$ -regularized ML estimator under the classical n > p setting, taking  $\theta_{\mathcal{S}}^{\natural}$  as the parameter to be estimated.

**Theorem 5.3.** We have

$$\left\|\check{\theta} - \theta^{\natural}\right\|_{\infty} \leq \frac{2}{\mu} \left( \left\| \left[ \nabla f(\theta^{\natural}) \right]_{\mathcal{S}} \right\|_{\infty} + \rho \right),$$

given that

$$\left\| \left[ \nabla f(\theta^{\natural}) \right]_{\mathcal{S}} \right\|_{\infty} + \rho \leq \frac{\mu}{2 \left\| A_{\mathcal{S}} \right\|_{\infty}}.$$

A short sketch of the proof of Theorem 5.3 is given in Section 6, where we make use of the fact that  $\overline{f} := \mathbb{E}[f]$  is self-concordant like.

Bounding  $\| [\nabla f(\theta^{\natural})]_{\mathcal{S}} \|_{\infty}$  by  $\| \nabla f(\theta^{\natural}) \|_{\infty}$ , we conclude that  $\hat{\theta} = \check{\theta}$  if the irrepresentability condition (Assumption 2) holds,  $\| \nabla f(\theta^{\natural}) \|_{\infty} \leq (\alpha/4)\rho$ , and

$$\rho \leq \min\left\{\frac{\alpha \mu^2}{2\lambda_{\max}(4+\alpha)^2 s \left\|A_{\mathcal{S}}\right\|_{\infty}}, \frac{2\mu}{(4+\alpha) \left\|A_{\mathcal{S}}\right\|_{\infty}}\right\}.$$

By applying Bernstein's inequality [14] to each element of  $\nabla f(\theta^{\natural})$  and the union bound, we show that  $\|\nabla f(\theta^{\natural})\|_{\infty}$  indeed concentrates around zero.

**Lemma 5.4.** For any t > 0 and any  $n > t^{-2}$ ,

$$\mathbb{P}\left\{\left\|\nabla f_n(\theta^{\natural})\right\|_{\infty} \geq t\right\} \leq 2p \exp\left[-\frac{\sqrt{n}t}{4\overline{\lambda}}\right].$$

where  $\overline{\lambda}^2 := \max\{1, \lambda_{\max}^2\}.$ 

*Remark.* Note that the concentration behavior is not sub-Gaussian like, which is assumed in [8].

Theorem 4.1 and Corollary 4.1 follow by combining the intermediate results above.

### 6. TECHNICAL SUPPLEMENTS

### 6.1. Self-concordant like functions

**Definition 6.1.** A function  $f : \text{dom}(f) \subseteq \mathbb{R}^p \to \mathbb{R}$  is selfconcordant like with parameter  $M \ge 0$  with respect to a norm  $\|\cdot\|$  on  $\mathbb{R}^p$  if dom(f) is open,  $f \in C^3(\text{dom}(f))$ , and

$$|D^3 f(x)[u, v, v]| \le M ||u|| D^2 f(x)[v, v]$$

for any  $x \in \text{dom}(f)$  and  $u, v \in \mathbb{R}^p$ .

*Remark.* The special case where  $\|\cdot\|$  is the  $\ell_2$ -norm is considered in [2, 24].

**Theorem 6.1.** Let  $f : \operatorname{dom}(f) \subseteq \mathbb{R}^p \to \mathbb{R}$  be a selfconcordant like function with parameter  $M \ge 0$  with respect to a norm  $\|\cdot\|$  on  $\mathbb{R}^p$ , and  $x, y \in \operatorname{dom}(f)$ . Define  $r := M \|y - x\|$  and the local norm

$$||y - x||_x := (D^2 f(x)[y - x, y - x])^{1/2}.$$

1. Bounds on the Hessian:

$$\exp(-r)\nabla^2 f(y) \le \nabla^2 f(x) \le \exp(r)\nabla^2 f(y).$$
 (3)

2. Bounds on the function value:

$$\omega_{*}(r) \|y - x\|_{x}^{2} \leq f(y) - f(x) - \langle \nabla f(x), y - x \rangle$$
  
$$\leq \omega(r) \|y - x\|_{x}^{2}, \qquad (4)$$

where  $\omega_*(r) := [\exp(-r) + r - 1]/r^2$  and  $\omega(r) := [\exp(r) - r - 1]/r^2$ .

*Proof.* Consider the function  $\psi_u(t) := D^2 f(y_t)[u, u]$  for any  $u \in \mathbb{R}^p$ , where  $y_t := x + t(y - x)$ . Then we have

$$|\psi'_{u}(t)| = |D^{3}f(y_{t})[y-x,u,u]| \le M ||y-x|| \psi_{u}(t),$$

and thus

$$|(\ln \psi_u(t))'| \le M ||y - x||.$$

We obtain (3) by integrating  $(\ln \psi_u(t))'$  over [0, 1]. We obtain (4) following the proof of Theorem 4.1.7 in [18].

# 6.2. Sketch of the Proof of Theorem 5.3

Define  $\overline{f}(\theta) := \mathbb{E}[f(\theta)]$ . Then we have  $\overline{f}(\theta) - f(\theta) = \langle \nabla f(\theta^{\natural}), \theta \rangle$ , and  $\theta^{\natural}$  minimizes  $\overline{f}(\theta)$ .

By definition

$$f(\theta) + \rho g(\theta) \le f(\theta^{\natural}) + \rho g(\theta^{\natural}),$$

which implies, by the triangle inequality,

$$\overline{f}(\hat{\theta}) - \overline{f}(\theta^{\natural}) \le \left( \left\| \left[ \nabla f(\theta^{\natural}) \right]_{\mathcal{S}} \right\|_{\infty} + \rho \right) \left\| \hat{\theta} - \theta^{\natural} \right\|_{1}.$$
(5)

It can be verified that for any  $x, u, v \in \mathbb{R}^p$  satisfying  $x_{S^c} = u_{S^c} = v_{S^c} \equiv 0$ ,

$$\begin{aligned} \left| D^{3}\overline{f}(x)[u,v,v] \right| &\leq \max\left\{ \left| (a_{i})_{\mathcal{S}}, u \right| \right\} D^{2}\overline{f}(x)[v,v] \\ &\leq \left\| A_{\mathcal{S}} \right\|_{\infty} \left\| u \right\|_{\infty} D^{2}\overline{f}(x)[v,v], \end{aligned}$$

and thus  $\overline{f}$ , being restricted on  $\mathbb{R}^{S}$ , is self-concordant like with parameter  $||A_{S}||_{\infty}$  with respect to the  $\ell_{\infty}$ -norm. By (4) and Assumption 1,

$$\overline{f}(\check{\theta}) - \overline{f}(\theta^{\natural}) \ge \frac{\mu}{M^2} \frac{\left\|\check{\theta} - \theta^{\natural}\right\|_1}{\left\|\check{\theta} - \theta^{\natural}\right\|_\infty} \left[\exp(-r) + r - 1\right], \quad (6)$$

where  $r := M \|\check{\theta} - \theta^{\natural}\|_{\infty}$  and  $M := \|A_{\mathcal{S}}\|_{\infty}$ . Combining (5) and (6), we obtain

$$\exp\left(-M\left\|\check{\theta}-\theta^{\natural}\right\|_{\infty}\right) + M\left\|\check{\theta}-\theta^{\natural}\right\|_{\infty} - 1$$
$$\leq \frac{M}{\mu}\left(\left\|\left[\nabla f(\theta^{\natural})\right]_{\mathcal{S}}\right\|_{\infty} + \rho\right)M\left\|\check{\theta}-\theta^{\natural}\right\|_{\infty}.$$
 (7)

Solving (7) directly, we obtain

$$M \left\| \check{\theta} - \theta^{\natural} \right\|_{\infty} \le W_0 \left[ -\frac{1}{1-a} \exp\left( -\frac{1}{1-a} \right) \right],$$

where  $W_0$  denotes the principal branch of the Lambert W function [6], and

$$a := \frac{M}{\mu} \left( \left\| \left[ \nabla f(\theta^{\natural}) \right]_{\mathcal{S}} \right\|_{\infty} + \rho \right)$$

We simplify the solution by Theorem 3.2 in [23], and the theorem follows.

#### 7. REFERENCES

- S. Arlot and A. Celisse, "A survey of cross-validation procedures for model selection," *Stat. Surv.*, vol. 4, pp. 40–79, 2010.
- [2] F. Bach, "Self-concordant analysis for logistic regression," *Electron. J. Stat.*, vol. 4, pp. 384–414, 2010.
- [3] P. Bickel, Y. Ritov, and A. B. Tsybakov, "Simultaneous analysis of Lasso and Dantzig selector," *Ann. Stat.*, vol. 37, no. 4, pp. 1705–1732, 2009.
- [4] P. J. Bickel, "Regularization in statistics," *Test*, vol. 15, no. 2, pp. 271–344, 2006.
- [5] A. C. Cameron and P. K. Trivedi, *Regression Analysis of Count Data*, 2nd ed. New York: Cambridge Univ. Press, 2013.
- [6] R. M. Corless, G. H. Gonnet, D. E. G. Hare, D. J. Jeffrey, and D. E. Knuth, "On the Lambert W function," Adv. Comput. Math., vol. 5, pp. 329–359, 1996.
- [7] B. Efron, "The estimation of prediction error: Covariance penalities and cross-validation," J. Am. Stat. Assoc., vol. 99, no. 467, pp. 619–632, Sep. 2004.
- [8] J. Fan and J. Lv, "Nonconcave penalized likelihood with NP-dimensionality," *IEEE Trans. Inf. Theory*, vol. 57, no. 8, pp. 5467–5484, Aug. 2011.
- [9] J. Fan and H. Peng, "Nonconcave penalized likelihood with a diverging number of parameters," *Ann. Stat.*, vol. 32, no. 3, pp. 928–961, 2004.
- [10] J. A. Fessler, "Statistical image reconstruction methods for transmission tomography," in *Handbook of Medical Imaging, Volume 2. Medical Image Processing and Analysis*, M. Sonka and J. M. Fitzpatrick, Eds. Bellingham: SPIE Press, 2000, ch. 1, pp. 1–70.
- [11] A. Gilbert and P. Indyk, "Sparse recovery using sparse matrices," *Proc. IEEE*, vol. 98, no. 6, pp. 937–947, Jun. 2010.
- [12] X. Jiang, G. Raskutti, and R. Willett, "Minimax optimal rates for Poisson inverse problems with physical constraints," 2014, arXiv:1403.6532v1 [math.ST].
- [13] Y.-H. Li, J. Scarlett, P. Ravikumar, and V. Cevher, "Sparsistency of ℓ<sub>1</sub>-regularized *M*-estimators," in 18th Inf. Conf. Artificial Intelligence and Statistics, 2015.
- [14] P. Massart, Concentration Inequalities and Model Selection. Berlin: Springer-Verl., 2007.
- [15] P. McCullagh and J. A. Nelder, *Generalized Linear Models*, 2nd ed. London: Chapman and Hall, 1989.

- [16] D. Motamedvaziri, M. H. Rohban, and V. Saligrama, "Sparse signal recovery under Poisson statistics," 2014, arXiv:1307.4666v2 [math.ST].
- [17] S. N. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu, "A unified framework for high-dimensional analysis of *M*-estimators with decomposable regularizers," *Stat. Sci.*, vol. 27, no. 4, pp. 538–557, 2012.
- [18] Y. Nesterov, Introductory Lectures on Convex Optimization. Boston, MA: Kluwer, 2004.
- [19] Y. Nesterov and A. Nemirovskii, *Interior-Point Polynomial Algorithms in Convex Programming*. Philadelphia, PA: SIAM, 1994.
- [20] M. Raginsky, R. M. Willett, Z. T. Harmany, and R. F. Marcia, "Compressed sensing performance bounds under Poisson noise," *IEEE Trans. Inf. Theory*, vol. 58, no. 8, pp. 3990–2010, Aug. 2010.
- [21] P. Ravikumar, M. J. Wainwright, and J. D. Lafferty, "High-dimensional Ising model selection using  $\ell_1$ regularized logistic regression," *Ann. Stat.*, vol. 38, no. 3, pp. 1287–1319, 2010.
- [22] P. Ravikumar, M. J. Wainwright, G. Raskutti, and B. Yu, "High-dimensional covariance estimation by minimizing *l*<sub>1</sub>-penalized log-determinant divergence," *Electron. J. Stat.*, vol. 5, pp. 935–980, 2011.
- [23] S. M. Stewart, "On certain inequalities involving the Lambert W function," J. Inequal. Pure Appl. Math., vol. 10, 2009.
- [24] Q. Tran-Dinh, Y.-H. Li, and V. Cevher, "Composite convex minimization involving self-concordant-like cost functions," 2015, arXiv:1502.01068 [math.OC].
- [25] S. van de Geer, "The deterministic Lasso," Seminar für Statistik, Eidgenössische Technische Hochschule, Research Report No. 140, 2007.
- [26] M. J. Wainwright, "Sharp thresholds for highdimensional and noisy sparsity recovery using  $\ell_1$ constrained quadratic programming (Lasso)," *IEEE Trans. Inf. Theory*, vol. 55, no. 5, pp. 2183–2202, May 2009.
- [27] E. Zeidler, Applied Functional Analysis: Main Principles and Their Applications. New York, NY: Springer-Verl., 1995.
- [28] P. Zhao and B. Yu, "On model selection consistency of Lasso," J. Mach. Learn. Res., vol. 7, pp. 2541–2563, 2006.